

METHOD AND APPARATUS FOR ENSURING MULTI-THREADED TRANSACTION  
ORDERING IN A STRONGLY ORDERED COMPUTER INTERCONNECT

CROSS-REFERENCE TO RELATED APPLICATIONS

5 This application is related to U.S. Application entitled "METHOD AND APPARATUS TO ALLOW DYNAMIC VARIATION OF ORDERING ENFORCEMENT BETWEEN TRANSACTIONS IN A STRONGLY ORDERED COMPUTER INTERCONNECT", filed on December 24, 2001, SN N/A, Attorney Docket No. 20206-28.

10 FIELD OF THE INVENTION

The present invention relates generally to transactions on a computer interconnect and, more specifically, to the ordering of read and write transactions on a computer bus.

15 DESCRIPTION OF THE RELATED ART

FIG. 1 shows the architecture of a typical computer system 8 in which a high-speed bus, such as the PCI bus 10, interconnects several I/O device adapters 12, 14. Each I/O device adapter 12, 14 is either an initiator or a target and the PCI bus (or PCI-X bus) serves to carry read and write transactions between the I/O device to which the adapter is connected. The CPU 16 for the computer system is connected to the bus 10 by means of a bridge device 18 which also provides a path between the CPU 16 and main memory 20. Another bridge device 22 connects a slower bus 24 to which devices, such as a printer adapter 26, and keyboard and mouse interfaces 28, are connected.

In one version of the PCI bus 10, an initiator (master) connects to a target (slave) via the bus to perform a transaction. FIG. 2 shows a typical PCI read transaction 40, a write transaction 42, and a retry request 44. Read transactions include an address phase 46, a command phase 48, one or more data phases 50a-d and attribute phases 52a-d. Each of the data phases 50a-d can be delayed by the target or initiator for a specific number of clocks in order to match the data transfer speed of the target to the initiator. Write transactions are similar, having an address phase 54, a command phase 56, one or more data phases 58a-d, and attribute phases 60a-d. The initiator or target can stall a data phase (via wait states) for up to seven clocks. (The target can

stall the start of the first data phase for up to 15 clocks). Before the initiator can connect to a target to perform a data transaction, the initiator must become the owner of the bus. This implies that the initiator must be the winner of an arbitration process.

In addition to data transactions, the PCI bus supports Delayed Transactions for reads and writes. A Delayed Transaction has two parts, the request part and the completion part. In the first part 44 in FIG. 2, the initiator performs an address phase 62 and command phase 64, and before the first data phase 66, the target responds with a disconnect 68, as shown in FIG. 2. The initiator interprets the target disconnect to be a retry request 44, which the initiator honors by ending the current transaction (FIG. 2), returning the bus to the idle state, re-arbitrating for ownership and re-initiating the transaction. If the initiator again receives a retry indication from the target, the initiator repeats the above sequence 44. Thus, the address phase 62, may be repeated several times (causing multiple re-arbitrations as well), until the target is ready to transfer data. In the completion part of the Delayed Transaction, the initiator performs an address phase and the target replies with a data transfer rather than a disconnect 68.

It is easily appreciated by one skilled in the art that the above-described operation of the PCI bus is exceedingly inefficient. Throughput on the PCI bus is lost for two reasons, the insertion of wait states and the use of the retry protocol.

Wait states cause a direct loss in throughput. Just one wait state inserted in each data phase is a 50% loss in throughput during the data burst. This means that for a 32 bit PCI bus clocked at 33 MHz, the throughput during the data phase is reduced to 66 Megabytes per second from 132 Megabytes per second. If the bus were clocked at 66 MHz, the throughput loss is even greater – a full 132 Megabytes per second of loss. For devices that can sustain transfer rates of about 1 Gigabyte per second, the bus is simply unworkable.

The Delayed Transaction protocol also causes a significant loss in throughput because bus cycles that could be used for data transfers are used to support a high-overhead protocol. Bus cycles are wasted when the target replies with a disconnect, when the initiator ends the current transaction, lets the bus go idle, re-arbitrates for the bus, and initiator then re-performs the address phase of the disconnected transaction. Thus, the cost of each retry is at least 6 clocks, 4 clocks to return the bus to the idle state, at least one clock for arbitration, and at least one more clock for an address phase. During these 6 clocks an entire 4 dword burst could have occurred.

In both of these cases throughput was lost because the target was not ready to respond. Clocking the bus faster to improve the throughput only causes more clock cycles to be lost due to wait states and the inefficient Delayed Transaction protocol.

An updated version of the PCI bus, PCI-X, was developed to address these and other 5 deficiencies. In the PCI-X specification, wait states are not permitted once data transfers have begun. A data burst, once started, must proceed at full speed on the bus. The read, write and split request transactions for PCI-X are shown for reference in FIG. 3. Each transaction type has an Address/Cmd phase, followed by an attribute phase and a response phase. After the response phase a data transfer ensues. Only the response and first data phase are extensible by adding a 10 limited number of wait states. After the first data phase, the remaining data phases must proceed at one bus clock per data phase.

Additionally, in the PCI-X specification, the inefficient Delayed Transactions have been replaced by Split Transactions, as shown in FIG. 3.

In a Split Transaction, a Requester initiates a transfer 75 by performing an address/cmd 15 phase 86, an attribute phase 88, a response phase 90, an unused data phase 92 and a surrender phase 94. Upon receiving a Split Response Request 96 from a Completer at the appropriate time in the transaction, the Requester removes itself from the bus, commits resources to the transaction and suspends the transaction until the Completer responds. This makes the bus 20 available for use to other Requesters and Completers in the interim. When the data transfer is ready to occur at the Completer, the Completer acts as an Initiator, obtaining the bus and performing a Split Completion transaction 71, which includes an address phase 70, an attribute phase 72, a response phase 74 and one or more data phases 76a-d during which the requested data is transferred to the Requester.

To make it easier to conform to the newer PCI-X specification and to improve the 25 performance of the older PCI protocol, it is best that both the Requester and Completer are implemented with read and write storage buffers so that when a write or read data burst is ready to occur, it can proceed at full bus speeds. Additionally, both the Requester and Completer are likely, in most implementations, to have Initiator and Target interfaces to carry out the Split Transaction protocol and each interface is required to be registered on both inputs and outputs.

However, the use of read and write buffers and Initiator and Target interfaces on the adapter increases the chances that PCI and PCI-X read/write ordering and deadlock avoidance rules may not be met.

PCI bus ordering rules require that if write data is posted to a write buffer (such as a posted-write buffer in a PCI-to-PCI or host/PCI bridge) the data must be flushed to its final destination (memory) before a read of that same data is allowed by the same or different bus master. Also, a bridge must perform all posted writes in the same order in which they were originally posted and is only permitted to post writes to regular memory targets.

On the PCI-X bus, there are more extensive read-write ordering rules when buffers are involved because of Split Transactions. For example, for bridges between a PCI-X bus and a host bus or between two PCI-X busses, there are three sets of rules, as set forth, in summary, below. The rules are set forth in more detail on pages 573-577 of PCI-X System Architecture, Tom Shanley, ISBN 0-201-72682-3, which is incorporated by reference into the present application.

Case I. A posted memory write transaction (PMW) is received in a bridge.

- (i) a subsequent split read request (SRR) or split write request (SWR) cannot be reordered to avoid returning incorrect read data (SRR) and to maintain write ordering (SWR);
- (ii) a subsequent split read completion (SRC) generally cannot be reordered to avoid returning incorrect read data;
- (iii) a subsequent split write completion may be permitted because the writes are in different directions; and
- (iv) a subsequent PMW generally cannot be performed until the first PMW is completed and PMWs must complete in the order received to maintain write ordering;

Case II - A split read request (SRR) or split write request (SWR) occurs at a bridge.

- (i) a subsequent split read request or split write request can be reordered;
- (ii) a subsequent split read completion, or split write completion or posted memory write must be allowed ahead of the SRR or SWR to avoid a deadlock;

Case III - A split read completion (SRC) or a split write completion (SWC) occurs at a bridge.

- (i) a subsequent split read request (SRR), split write request (SWR), split read completion (SRC) or split write completion (SWC) can be reordered;

(ii) a posted memory write must go ahead of the SRC or SWC to avoid a deadlock.

Though there are some exceptions to these rules if a relaxed ordering (RO) bit is set in a transaction, the rules set forth the ordering of reads and writes so as to guarantee the proper 5 operation of system software (Case I) and the avoidance of deadlocks (Cases II and III) for a bridge between the host bus and a PCI-X bus or between two PCI-X busses.

The use of read and write buffers and both an Initiator and Target interface on adapter units connected to the PCI bus along with the split transaction protocol creates a need to maintain certain ordering of read and write transactions on the adapter units.

10 Currently, one method of dealing with the ordering problem is to control all of the read-write transaction activity from a single thread, thereby serializing all of the transactions from a single point of control. While this may assure that the ordering problem is correctly addressed, the single thread approach is performance limiting both to the adapter and the system.

15 Therefore, there is a need to address the ordering problem on an adapter that has read and write buffers, Initiator and Target interfaces that connect to a strongly-ordered bus, such as the PCI or PCI-X bus, and a split transaction protocol, without using a single thread to serialize all of the activity of the adapter. The present invention is directed towards such a need.

## BRIEF SUMMARY OF THE INVENTION

20 The present invention provides an apparatus and method for meeting the ordering rules of PCI and PCI-X on an adapter unit connected to the PCI or PCI-X bus. The adapter unit has both an Initiator and Target interface for handling delayed or split transactions. In addition, the adapter unit is capable of multi-threaded operation, supporting more than one read or write transaction thread.

25 In accordance with one version of the present invention, buffer sharing is used to enforce ordering between split or delayed read completions (SRC, supplied by the target interface) and a write transaction (supplied by the initiator interface and posted in the adapter unit's local write buffer) and between two successive write transactions. The SRCs are queued in the same buffer as the write transactions so that they occur on the PCI bus in the order in which they originally 30 occurred. This meets the Case I (ii) rule.

Ordering is enforced between initiator read transaction requests and posted initiator write transactions by write contingency flags. Whenever the initiator interface of the adapter unit issues a read transaction request, write contingency flags are set for the read transaction request, one for each pending initiator write that was accepted prior to the requested initiator read

5 transaction request but not performed on the PCI bus. When a write transaction is performed on the bus, the contingency flag corresponding to the write transaction is cleared in all pending initiator read transaction requests and when all the write contingency flags have been cleared for a particular initiator read transaction request, the initiator read transaction request is then permitted on the bus. This meets the Case I(i) rule.

10 In another version of the invention, read transaction requests are either ordered or non-ordered. Ordered read transaction requests have their write contingency flags set, whereas non-ordered read transaction requests do not. Whenever an ordered read transaction request is accepted, write contingency flags are set in the read transaction context for each pending write, i.e., writes that were accepted prior to the read and not yet performed on the PCI bus. Whenever a pending write is performed on the PCI bus, the contingency bits in the read context are cleared. For non-ordered read transaction requests, the write contingency flags are not set. A read transaction request is not presented on the PCI bus until all write contingencies are cleared for that read transaction request.

15

20 In yet another version of the invention, read transaction requests are selectively ordered, i.e., the read transaction requests are dependent on the retiring of selected ones of the pending write transactions. This selectively ordering is based on the properties of either the read transaction request, or one or more pending write transactions, or both. Properties of these transactions for the purpose of selectively ordering the read and write transactions include, but are not limited to, the address or address range of the transaction, and transaction attributes, such 25 as the requester function number, the requester device number and the relaxed ordering attribute bit.

30 One method in accordance with the present invention is a method for enforcing ordering between read and write transactions for an adapter unit configured to connect to a bus, the adapter unit having at least one read and at least one write buffer. The method includes the steps of enqueueing, in the write buffer, an initiator write transaction to be performed on the bus,

receiving, in the read buffer, an initiator read transaction request to be performed on the bus, where the performance, on the bus, of the initiator read transaction request is required to follow the performance, on the bus, of the initiator write transaction, and clearing a write contingency flag in a register associated with the read transaction request. If, after receiving the read  
5 transaction request, the performance of the initiator write transaction has not started on the bus, the method includes setting the write contingency flag in the register to indicate the dependence of the initiator read transaction request on the initiator write transaction in the write buffer, performing the initiator write transaction on the bus, and after the performance of the write transaction is completed, clearing the write contingency flag in the register. If the write  
10 contingency flag is not set, then the read transaction request is performed on the bus.

Additionally, the method includes enqueueing a target read completion transaction in the write buffer, where the target read completion transaction has a position in the write buffer either before or after the initiator write transaction based on whether the target read completion originally occurred before or after the initiator write transaction, and performing on the bus the transactions in the write buffer in the order in which the transactions were enqueued.  
15

One advantage of the present invention is that the ordering rules of the PCI or PCI-X bus are not exported to any of the data sources or sinks with which the adapter unit of the present invention interoperates.

Another advantage of the present invention is that the flexibility of using ordered or non-  
20 ordered read transaction requests permits higher bandwidth and eliminates the chance of a deadlock condition. Transactions that are independent of each other need not be burdened with the ordering rules, whereas transactions that are inter-dependent are burdened with the ordering rules so that correct software operation is assured.

Yet another advantage is that a multi-threaded adapter unit can be employed on the PCI  
25 bus. This type of adapter unit permits more than one outstanding split read or write transaction on the PCI bus without risk of an ordering violation or deadlock.

#### BRIEF DESCRIPTION OF THE DRAWINGS

These and other features, aspects and advantages of the present invention will become better understood with regard to the following description, appended claims, and accompanying drawings where:

FIG. 1 shows the architecture of a typical computer system in which a high-speed bus, such as the PCI bus, interconnects several I/O device adapters;

FIG. 2 shows a typical PCI read, write and retry operation;

FIG. 3 shows a typical PCI-X read, write and split response operation;

FIG. 4 shows the configuration of the adapter unit interface to the PCI bus;

FIG. 5 shows a flow chart for handling initiator read and write transactions and target read completions; and

FIG. 6 show a flow chart for operating the write buffer in an adapter unit.

#### DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENTS

FIG. 4 shows the configuration of the adapter unit interface 100 to the PCI bus. The adapter unit 100 includes a PCI buffer manager section 102 that has two interfaces 104, 106. A first interface 104 is an internal interface to the functional sections of an application specific device. A second interface 106 is an interface to a PCI bus protocol block 108, which connects to the PCI bus and is responsible for carrying out the PCI specific protocols, such as dword commands and burst commands. The latter include memory read, memory write, memory read block, memory write block and delayed and split completions.

Included in the PCI buffer management section 102 are a plurality of read buffers 110 and write buffers 112, a set of write context registers 114 for tracking the contents of the write buffers, a set of read context registers 116 for tracking write contingency flags for the read buffers 110, and a plurality of multiplexors 118a-e or data selectors for routing data. Each write context register 114 is associated with one entry of the M write buffers 112 and contains a bit (one of sr bits 120) that indicates whether the associated entry is an initiator write transaction or a target read completion transaction. Each read context register 116 is associated with one entry of the read buffers 110 and contains M wc bits 122, each of which, if set, indicates the presence of a single pending write in the write buffer that is ahead of the initiator read transaction request in the read buffer. As is shown in FIG. 4, there are as many wc bits (M) in the read context

register as there are write buffers (M). This arrangement permits a read transaction request in the read buffer to be dependent on selected ones of the write transactions in the write buffer. The decision of whether to set a wc bit to create an ordering between a write transaction and a read transaction request is based on the properties of the read and write transactions involved. Such 5 properties can include the address or address ranges of the transactions and the attributes of the transactions such as the requester function or device number in the transaction.

A number of multiplexors (or data selectors) are provided including multiplexor 118a for providing read buffer data to one or more read masters in the application specific section of the adapter, multiplexer 118b for selecting a write master to provide data to the write buffers, 10 multiplexer 118d for selecting a read slave to provide data to the write buffers and a multiplexer 118c for selecting whether a read master or a write master has access to the write buffers. Finally, a multiplexer 118e is provided for enabling write buffer data to either an initiator interface part or a target interface part of interface 106 via register 124, depending on the sr flags 120 stored in the write context registers.

15 FIG. 5 shows a flow chart for handling initiator read and write transactions and target read completions. When a read transaction request appears on the internal interface, as determined in step 140, all of the write contingency flags for the read transaction request are cleared, in step 142, and if ordering is enforced, as determined in step 144, then a test is made, in step 146, to determine whether there are any pending initiator write transactions. These are write transactions in the write buffers that have been accepted by the write buffer before the 20 appearance of the read transaction request but whose performance on the PCI bus has not started. If there are pending initiator write transactions, then a write contingency flag for each pending write transaction on which the read transaction request must depend is set, in step 148, in the read context register for the read transaction request that just appeared. The read transaction 25 request then waits, in step 150, for all of the write contingency flags that were set for the read transaction request to be cleared, before being presented to the initiator interface and thus to the PCI bus, in step 152. The contingency flags thus assure proper ordering between initiator read transaction requests and initiator write transactions. If ordering is not enforced, as determined in step 144, or there are no pending initiator write transactions for the read transaction request, as 30 determined in step 146, then the read transaction request proceeds without delay, in step 152. If

ordering is not enforced, improvements in throughput of the adapter unit are possible, because read transaction requests are not delayed by any pending write transactions. Selective ordering is also possible, whereby the read transaction request may be ordered with respect to some pending write transactions and not others. In this case, only those write contingency flags pertaining to 5 pending write transactions on which the read transaction request must wait are set.

If a write transaction is presented on the internal interface, as determined in step 140, the write transaction is enqueued, in step 154, in one of the available write buffers. If a read transaction completion appears on the internal interface, as determined in step 140, it too is enqueued, in step 154, in one of the available write buffers and the sr bit is set, in step 158, as 10 determined in step 156. This transaction is a delayed or split completion to a read transaction request that was issued by a source external to the adapter unit, which signaled a delayed or split response to the read request. The sr bit in the write context registers keeps track of whether the entry in the write buffer is a write transaction or a delayed or split completion transaction. The queuing of initiator write transactions and target read completions assures that these transactions 15 are performed in the order in which they are enqueued in the write buffer.

FIG. 6 show a flow chart for operating the write buffers 112 in an adapter unit. In this process, if there are any transactions in the write buffers 112, as determined in step 170, the transaction at the head of the queue is performed on the PCI bus by presenting it to the PCI Bus Protocol block 108. If a write transaction is the transaction performed in step 172, as determined in step 174, then the write contingency flags for that write transaction are cleared, in step 176, in each read context register in which they are set, i.e., for all pending initiator read transaction requests that are subsequent to the write transaction that just completed. After all write transactions, on which a given read transaction request is contingent, have been performed on the PCI bus and their write contingency flags cleared, in step 150 in FIG. 5, the read transaction 20 25 request has no contingencies and can proceed to the PCI bus.

Although the present invention has been described in considerable detail with reference to certain preferred versions thereof, other versions are possible. Therefore, the spirit and scope of the appended claims should not be limited to the description of the preferred versions contained herein.